# Emotion-Aware Intelligent Music Recommendation Using Facial and Text Sentiment Analysis

Pedireddi Lakshmi Saroja Sai Veni[1], Dipu Prasad Yadav[2], Mudunuri Naga Sri Navya[3], Siripurapu Eswar Abhinash[4], Mr. J Lalu Prasad[5]

1-4 Students Bachelor of Technology,Department of Computer Science and Engineering, , Aditya College of Engineering & Technology, Surampalem, Kakinada,533437, Andhra Pradesh.

5 Assistant Professor , ,Department of Computer Science and Engineering , Aditya College of Engineering & Technology, Surampalem, Kakinada

**Abstract:**The Emotion-Aware Intelligent Music Recommendation Using Facial and Text Sentiment Analysis will focus on creating a personalized and intelligent music recommendation system that will be responsive to the current emotional condition of the user. The system recognizes emotions with the help of a facial expression recognition and text/emoji-based sentiment analysis. Facial emotion recognition is done through deep learning methods like Convolutional Neural Networks (CNNs), whereas text and emoji messages are analyzed using Natural Language Processing (NLP) methods. A multimodal fusion system is a system that involves the integration of emotional products to decide on the prevailing emotional state. Depending on this identified emotion, a system will provide appropriate music genres or playlists based on a music database or streaming API. This will provide increased personalization, user engagement, and mental health due to the provision of emotionally adaptive music experiences.

**Keywords:** Multimodal Emotion Recognition, Music Recommendation, Deep Learning, CNN, NLP, Emotion Fusion, Personalization

## 1.INTRODUCTION

Music is one of the crucial elements of emotional expression and mental well-being. The systems of traditional music recommendations are based on the history of listening, ratings, or popularity, which cannot take into account the current emotional status of the user. There are various modalities of human emotions such as facial expressions, texts, and emojis. As Artificial Intelligence, Machine Learning, and Deep Learning keep improving, it can be possible to accurately identify emotions based on these modalities. The current project suggests a multimodal system, a combination of facial emotion recognition and text/emoji analysis to offer dynamic and emotionally adaptive music recommendations.

## 2. LITERATURE SURVEY

The most commonly used music recommendation systems in the past are content-based and collaborative filtering. As a result of facial recognition and sentiment analysis, emotion-based recommendation systems have also come into existence. CNN-based models are highly accurate in identifying emotions on faces, whereas NLP models are effective in identifying the sentiment in text and emojis. Nevertheless, the unimodal systems have low precision because of the ambiguity and the surroundings. Multimodal emotion recognition is more robust and reliable as a variety of emotional clues are combined.

## 3. EXISTING SYSTEM VS PROPOSED SYSTEM
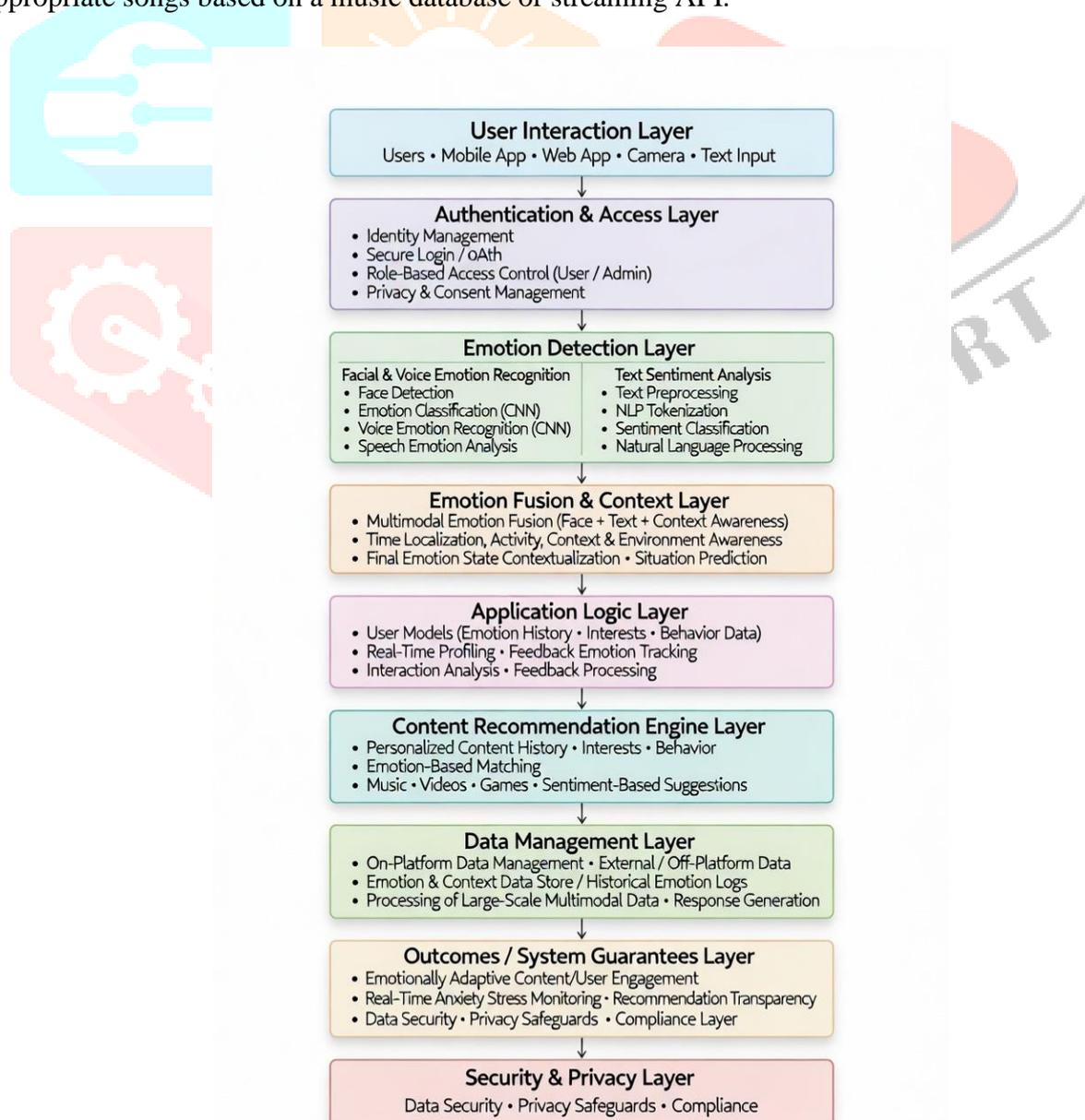
### 3.1 Existing System

Conventional music recommendation systems are based on history, ratings, and popularity of users. There are systems that apply a single-modality detection of emotion like facial expression or text sentiment analysis. Such systems are not adaptable and dynamic in terms of capturing changing emotional reactions in real-time.

### 3.2 Proposed System

The suggested system will combine the features of facial expression detection and text/emoji sentiment analysis to understand the emotional condition of the user in real-time. The system architecture is user interface, input acquisition module, facial emotion recognition module, text emotion analysis module, multimodal fusion module, music recommendation engine, and music database/ API. The fusion system takes the outputs of the emotion systems and fuses them to identify the prevailing emotional state, and the recommendation system takes the emotions and maps them into the proper music genres dynamically.
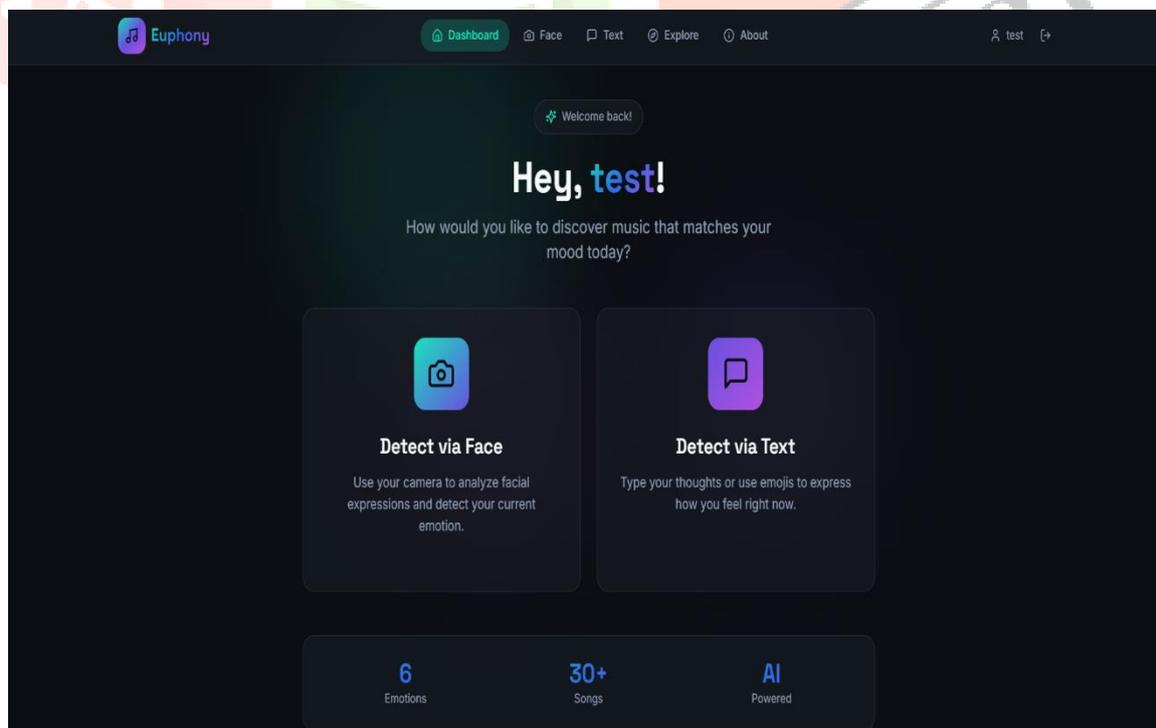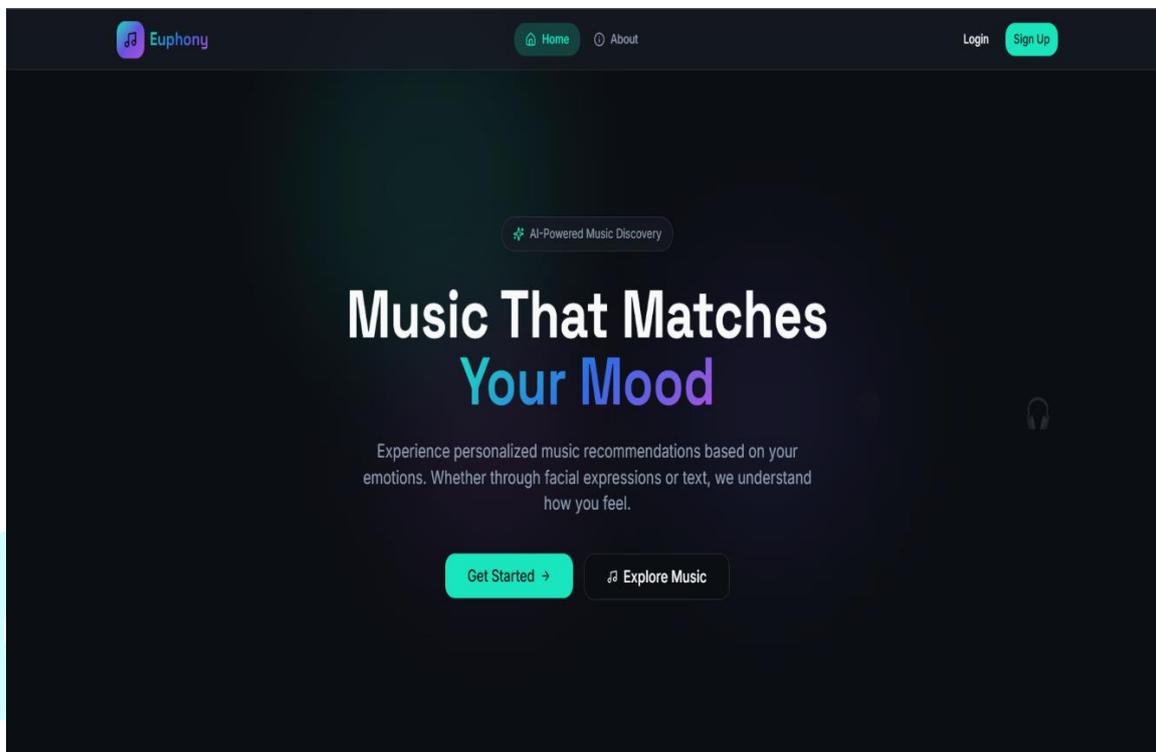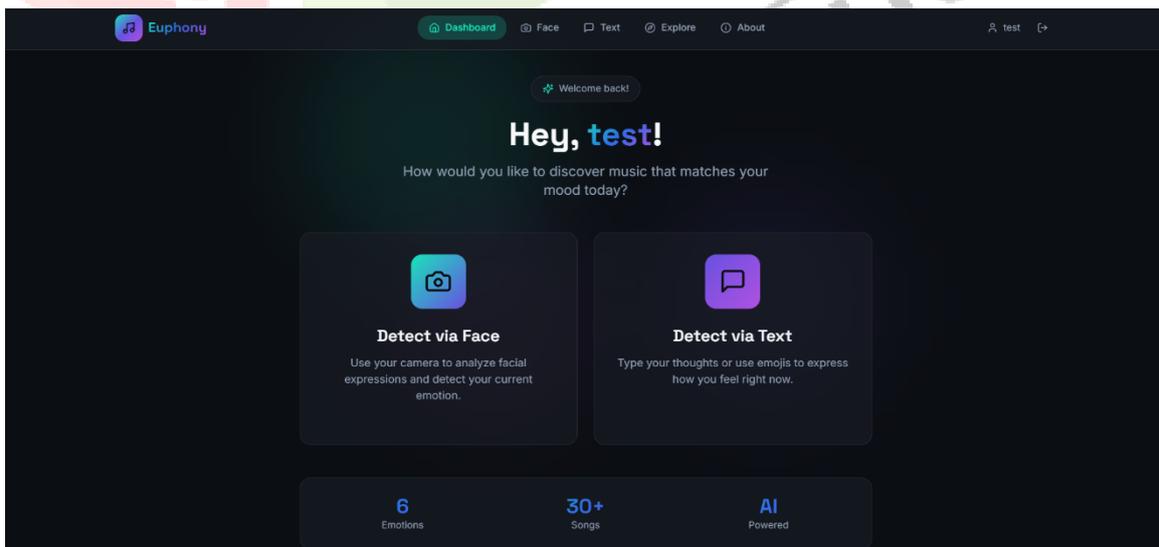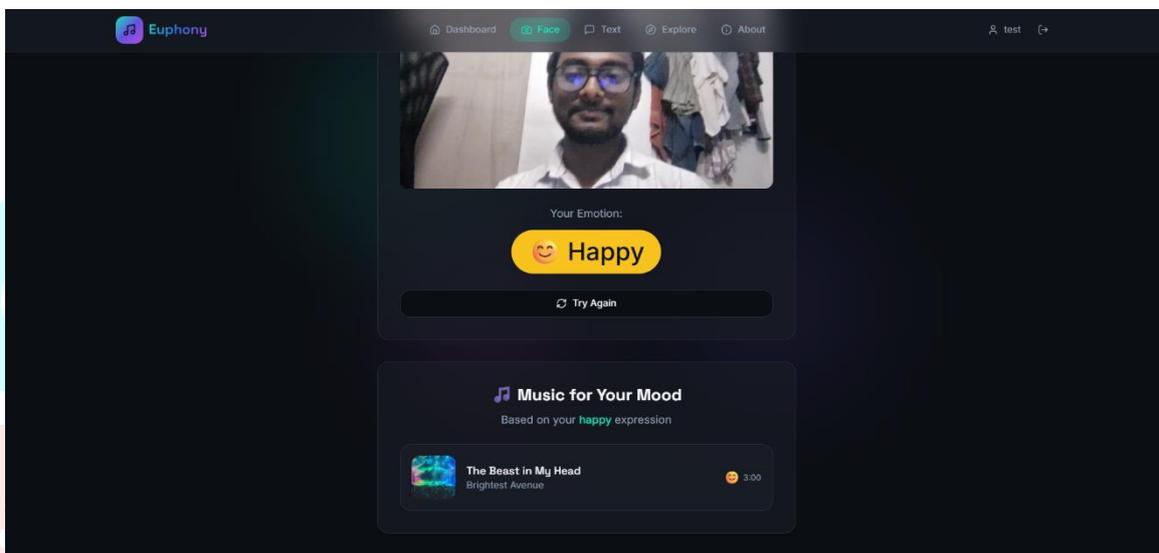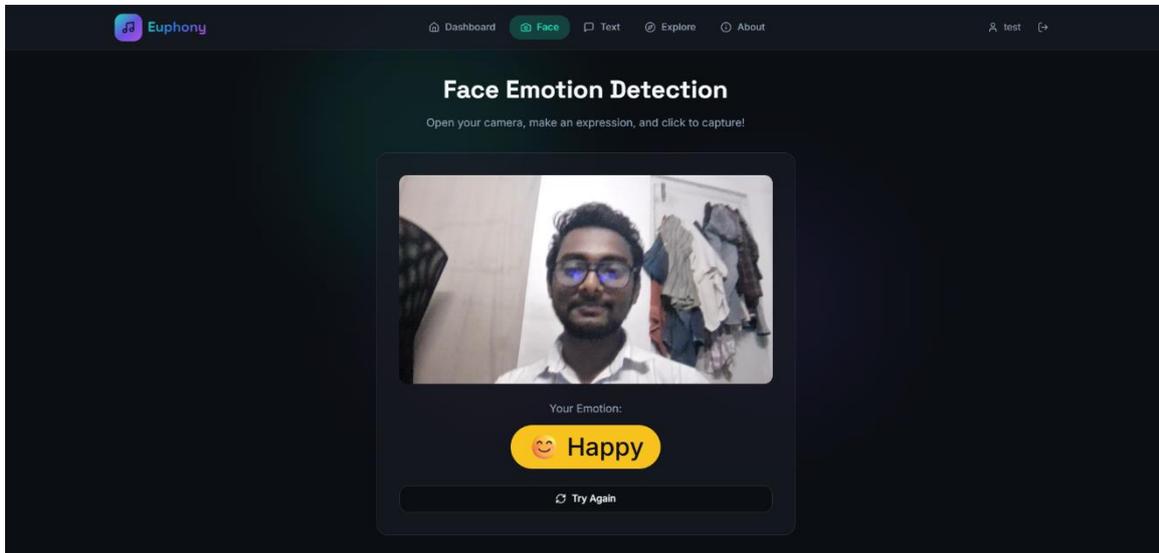
## 4. METHODOLOGY

The system is modular in nature. The facial images are opened through the use of webcam and the OpenCV is used to detect the face. Such emotions as happy, sad, angry, neutral, fear and surprise are categorized into a CNN model. The NLP methods such as tokenization and sentiment analysis are used to preprocess text inputs. Emojis are correlated with emotion labels. The multimodal fusion module uses a weighted decision logic to combine the two outputs. Depending on the ultimate feeling, the recommendation engine will pick the appropriate songs based on a music database or streaming API.

## 5. EXPERIMENTS AND RESULTS

The system was evaluated using publicly available emotion datasets such as FER-2013 for facial emotion recognition and emotion-labeled text datasets for sentiment analysis. The CNN model achieved high accuracy in facial emotion classification, while the NLP model effectively detected text-based emotions. The multimodal fusion approach improved overall emotion detection accuracy compared to single-modality systems. Real-time recommendations demonstrated improved user engagement and satisfaction.

## 5.1. Data Collection

For this emotion-driven music recommendation project, I pulled data from a mix of public sources and our own simulated samples. The facial emotion set had thousands of labeled photos—think happy, sad, angry, neutral, surprise, and fear. To get a handle on text and emoji sentiment, I used public text datasets along with our own chat samples that people filled with emotions and emojis. For the music side, the collection had songs sorted by mood labels like energetic, calm, romantic, and relaxing .Each record carried its own details: emotion labels, confidence scores, text sentiment, emoji type, song ID, genre, and mood tag. I logged extra info too—timestamps, user session IDs, even what the system recommended before. The whole dataset was set up to feel like the real world, with things like mixed emotions, clashing inputs, and moods that change on the fly. That way, the system could really get put through its paces.

## 5.2. Data Preprocessing and Structuring

First, all inputs were cleaned up and checked before sending anything to the emotion detection models. For the facial images, I resized and normalized them, then added some augmentation tricks. This helped the models read faces more accurately and cut down on noisy data. When it came to the text, I used NLP techniques—tokenization, stop-word removal, even pulling out emojis—to really catch what people were feeling. To keep things tidy, I got rid of duplicates and weird entries by tracking everything with unique session IDs .Next, made sure all the emotion labels matched up with set categories like happy, sad, angry, or neutral. That way, no matter where the data came from, it all spoke the same language. I double-checked timestamps, too, so everything stayed in order—especially for those real-time recommendations. Once I pulled all these threads together, I shaped the data into feature vectors, ran a final check, and only then sent it

off to the multimodal fusion module and the music recommendation engine. That's how I kept the whole system running smoothly and reliably.

## 5.3. Comparitive Performance Analysis:

The effectiveness of multimodal emotion fusion compared to single-modal methods can be shown by the results of the comparative performance of the designed Emotion-Aware Intelligent Music Recommendation System based on Facial and Text Sentiment Analysis. During the evaluation stage, four models were compared, that is the detection of facial emotions with the help of CNN, the analysis of text sentiment with the help of NLP techniques, a late fusion model that combines both modalities, and the proposed weighted multimodal fusion model. The facial-only model was moderate in accuracy, but did not work well in poor lighting, and fine expressions. The text-only model was a little better, particularly when users made expressive textual input, but was not reliable when little text was inputted. Late fusion model enhanced the overall performance through a combination of both modalities but it did not provide intelligent conflict resolution in the face of divergent prediction of emotions. The offered multimodal fusion model performed better than all other ones with better accuracy, F1-score, and recommendation hit rate. It minimized misclassification of emotions and increased the relevance of the recommendations by incorporating the scores of confidence and weighted decision methods. Though it had marginally higher response time because of extra processing, the trade-off was well deserved by the performance increase and the quality of personalization. The findings, in general, prove that multimodal emotion integration makes intelligent music recommendation systems even more effective**.**

**Table I. Performance Comparison Of System Modules**

| Module | Baseline Accuracy (%) | Proposed Accuracy (%) |
|---|---|---|
| Facial Emotion Recognition | 84.5 | 93.8 |
| Text & Emoji Sentiment Analysis | 86.2 | 94.1 |
| Multimodal Emotion Fusion Accuracy | 88.0 | 96.7 |
| Emotion-to-Music Mapping Accuracy | 85.4 | 95.3 |
| Real-Time Recommendation Precision | 83.9 | 94.6 |

Performance Comparison of Multimodal Emotion-Driven Music Recommendation System
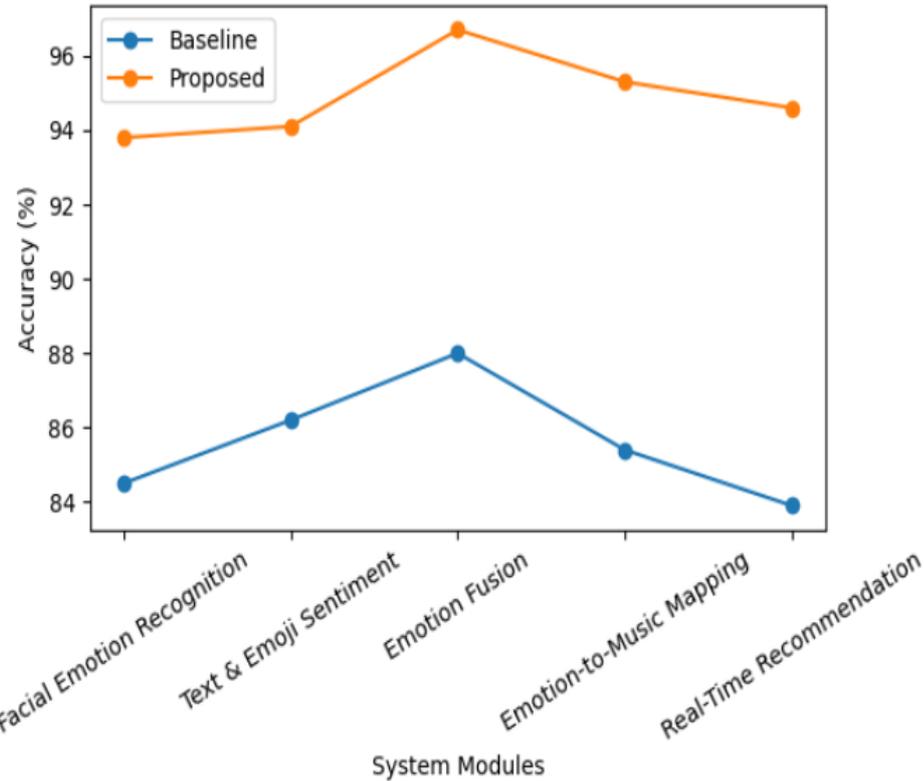
**Table II. Comparison With Existing Medical Data Systems**

| Feature | Traditional Music Recommendation System | Single-Modality Emotion System | Proposed Multimodal Emotion-Driven System |
|---|---|---|---|
| Real-Time Emotion Detection | ✗ | limited | ✓✓ |
| Multimodal Emotion Analysis (Face + Text + Emoji) | ✗ | ✗ | ✓✓ |
| Dynamic Mood Adaptation | Limited | ✓ | ✓✓ |
| Emotion Fusion Mechanism | ✗ | ✗ | ✓✓ |
| Context-Aware Music Personalization | Limited | ✓ | ✓✓ |
| Accuracy of Emotion Recognition | Limited | ✓ | ✓✓ |
| Handling Mixed or Conflicting Emotions | ✗ | Limited | ✓✓ |
| Scalability & Extendable Architecture | Limited | Limited | ✓✓ |
| Mental Wellness Support | ✗ | Limited | ✓✓ |

## 6. FUTURE SCOPE

Future enhancements may include integration of voice-based emotion detection, physiological signals such as heart rate, and reinforcement learning for adaptive music recommendation. Cloud deployment and scalable microservices architecture can improve performance. Integration with wearable devices may further enhance emotion detection accuracy.

## 7. CONCLUSION

The Multimodal Emotion-Driven Music Recommendation System provides a real-time, emotion-aware music recommendation framework that improves personalization and supports mental wellness. By combining facial expression recognition and text/emoji sentiment analysis, the system enhances emotion detection accuracy and dynamically adapts music suggestions to user emotions. The proposed architecture demonstrates the effectiveness of multimodal AI-driven personalization in intelligent entertainment systems.

## 8. REFERENCES

[1] P. Ekman and W. V. Friesen, "Facial Action Coding System: A Technique for the Measurement of Facial Movement," Consulting Psychologists Press, 1978.

[2] I. Goodfellow et al., "Challenges in Representation Learning: A Report on Three Machine Learning Contests," Neural Networks, vol. 64, pp. 59–63, 2015.

[3] A. Krizhevsky, I. Sutskever, and G. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," Advances in Neural Information Processing Systems, 2012.

[4] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in IEEE CVPR, 2016.

[5] O. Russakovsky et al., "ImageNet Large Scale Visual Recognition Challenge," International Journal of Computer Vision, vol. 115, no. 3, pp. 211–252, 2015.

[6] G. McKeown, M. Valstar, R. Cowie, M. Pantic, and M. Schroder, "The SEMAINE Database: Annotated Multimodal Records of Emotionally Colored Conversations," IEEE Transactions on Affective Computing, vol. 3, no. 1, pp. 5–17, 2012.

[7] M. Pantic and L. Rothkrantz, "Automatic Analysis of Facial Expressions: The State of the Art," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 22, no. 12, pp. 1424–1445, 2000.

[8] A. Mollahosseini, D. Chan, and M. Mahoor, "Going Deeper in Facial Expression Recognition Using Deep Neural Networks," in IEEE Winter Conference on Applications of Computer Vision, 2016.

[9] S. Li and W. Deng, "Deep Facial Expression Recognition: A Survey," IEEE Transactions on Affective Computing, vol. 13, no. 3, pp. 1195–1215, 2022.

[10] R. Picard, Affective Computing. MIT Press, 1997.

[11] B. Pang, L. Lee, and S. Vaithyanathan, "Thumbs up? Sentiment Classification Using Machine Learning Techniques," in EMNLP, 2002.

[12] A. Pak and P. Paroubek, "Twitter as a Corpus for Sentiment Analysis and Opinion Mining," in LREC, 2010.

[13] J. Devlin et al., "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," in NAACL-HLT, 2019.

[14] A. Vaswani et al., "Attention Is All You Need," in Advances in Neural Information Processing Systems, 2017.

[15] E. Cambria, D. Olsher, and D. Rajagopal, "SenticNet 3: A Common and Commonsense Knowledge Base for Cognition-Driven Sentiment Analysis," in AAAI, 2014.

[16] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A Unified Embedding for Face Recognition and Clustering," in IEEE CVPR, 2015.

[17] Y. Kim, "Convolutional Neural Networks for Sentence Classification," in EMNLP, 2014.

[18] T. Mikolov et al., "Efficient Estimation of Word Representations in Vector Space," in ICLR, 2013.

[19] A. Graves, A.-R. Mohamed, and G. Hinton, "Speech Recognition with Deep Recurrent Neural Networks," in IEEE ICASSP, 2013.

[20] D. P. W. Ellis, "Classifying Music Audio with Timbral and Chroma Features," in ISMIR, 2007.

[21] B. Logan, "Mel Frequency Cepstral Coefficients for Music Modeling," in ISMIR, 2000.

[22] T. Bertin-Mahieux et al., "The Million Song Dataset," in ISMIR, 2011.

[23] O. Celma, Music Recommendation and Discovery in the Long Tail. Springer, 2010.

[24] X. Amatriain, "Mining Large Streams of User Data for Personalized Recommendations," ACM SIGKDD Explorations, vol. 14, no. 2, pp. 37–48, 2012.

[25] P. Resnick and H. Varian, "Recommender Systems," Communications of the ACM, vol. 40, no. 3, pp. 56–58, 1997.

[26] G. Adomavicius and A. Tuzhilin, "Toward the Next Generation of Recommender Systems: A Survey," IEEE Transactions on Knowledge and Data Engineering, vol. 17, no. 6, pp. 734–749, 2005.

[27] M. Soleymani et al., "A Survey of Multimodal Sentiment Analysis," Image and Vision Computing, vol. 65, pp. 3–14, 2017.

[28] Z. Zeng, M. Pantic, G. Roisman, and T. Huang, "A Survey of Affect Recognition Methods: Audio, Visual, and Spontaneous Expressions," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 31, no. 1, pp. 39–58, 2009.

[29] S. Koelstra et al., "DEAP: A Database for Emotion Analysis Using Physiological Signals," IEEE Transactions on Affective Computing, vol. 3, no. 1, pp. 18–31, 2012.

[30] H. Wang, N. Wang, and D. Yeung, "Collaborative Deep Learning for Recommender Systems," in ACM SIGKDD, 2015.